

# LECTURE 4 : GAUSSIAN DISTRIBUTION

## [r] Motivation

Both MLE & MAP give us the prediction of  $\bar{w}^*$  in the training stage :  $\bar{w}^* = \underset{\{\bar{w}\}}{\operatorname{argmax}} p(\bar{w} | \bar{x}, \bar{t})$ .

This  $\bar{w}^*$  is later used in the testing stage :

$$t_{\text{new}} = \bar{w}^* \bar{x}_{\text{new}} \quad \text{or} \quad t_{\text{new}} = \bar{w}^* \phi(\bar{x}_{\text{new}})$$

Q: So, what's wrong with MLE or MAP?

A1: Both MLE & MAP provide the prediction ( $\bar{w}^*$ ). but they don't provide any information on the probability or uncertainty information. Namely, how certain the predicted  $\bar{w}^*$  is.

A2: This probability information is important in decision making in general, but also critical when we need to compare two or more prediction results. Of course, we will choose or be more confident to use the predicted result with the highest probability or certainty.



Full Bayesian Inference resolves the problem by using :

$$p(x|d) = \frac{p(d|x) p(x)}{p(d)}$$

where :  $p(d) = \int p(d, x) dx \approx \sum_x p(d, x)$

$$= \int p(d|x) p(x) dx \approx \sum_x p(d|x) p(x)$$

Q : To obtain  $p(x|d)$  requires  $p(d)$ . Then, how to obtain  $p(d)$  ?

A : 
$$p(d) = \int p(d|x) p(x) dx \approx \sum_x p(d|x) p(x)$$

There are many possible functions for  $p(d|x) p(x)$

In discrete, it can be computationally expensive.



Many of them are difficult to be integrated (analytically intractable).

e.g. :  $\int \exp(x^3 + x^2) dx$

Q : Why is  $\sum_x p(d|x) p(x)$  computationally prohibitively expensive?

A : Let's make  $p(d) = \sum_{\bar{x}} p(d, \bar{x})$ , where  $\bar{x}$  is a vector.

If  $\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , then :

$$p(d) = \sum_{x_1} \sum_{x_2} p(d, \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}) = \sum_{x_1} \sum_{x_2} p(d, x_1, x_2)$$

In programming it means :

```

for all candidates of ( $x_1$ ) do
  for all candidates of ( $x_2$ ) do
     $p(d) = p(d) + p(d, x_1, x_2)$ 
  end;
end;
  }
```

$\} (\# \text{candidates})^2$

Imagine if  $\bar{x}$  is a vector with length of 100 & the number of candidates is 50, then the complexity is  $(50)^{100}$



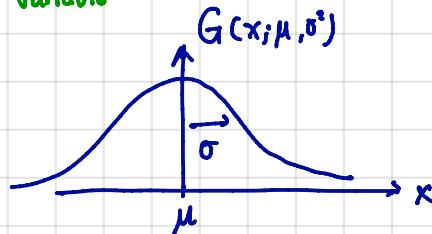
This is the reason why it is prohibitively expensive!

## [2] Univariate & Multivariate Gaussian

Univariate :

$$G(x; \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

a scalar variable

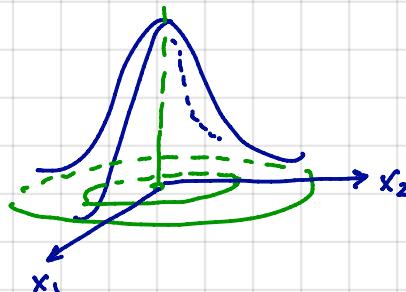


Multivariate :

$$G(\bar{x}; \bar{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (\bar{x} - \bar{\mu})^\top \Sigma^{-1} (\bar{x} - \bar{\mu})\right)$$

$\downarrow$   
a vector

covariance matrix



Covariance Matrix,  $\Sigma$  :

Definition of covariance :  $\text{cov}(a, b) = E[(a - E(a))(b - E(b))]$

$\downarrow$

Expectation

Covariance Matrix

$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

$\rightarrow$  the covariance of  $x_i$  &  $x_j$  :

$$\Sigma_{ij} = \text{cov}(x_i, x_j) = E[(x_i - \mu_i)(x_j - \mu_j)]$$

where :  $\mu_i = E[x_i]$

the average of  $(x_i^1, x_i^2, \dots, x_i^N)$

$\downarrow$

$$\text{Hence : } \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \dots & \Sigma_{1D} \\ \Sigma_{21} & \ddots & \ddots & \vdots \\ \vdots & & & \\ \Sigma_{D1} & \dots & \dots & \Sigma_{DD} \end{bmatrix}$$

$$\text{cov}(a, b) = E[ab] - E[a]E[b]$$

Question : What is the geometrical meaning of a covariance matrix?

### [3] Covariance Matrix : Geometric Meaning

#4

Let  $\bar{x} \in \mathbb{R}^2$  :

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}_{2 \times 2}$$

→ a covariance matrix is always a symmetric matrix:

$$\begin{aligned} \Sigma_{21} &= \Sigma_{12} \\ \downarrow & \\ \Sigma^T &= \Sigma \end{aligned}$$

Assuming  $\bar{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ :

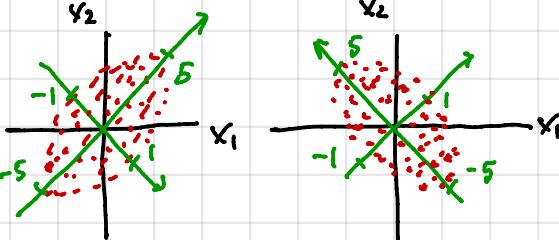
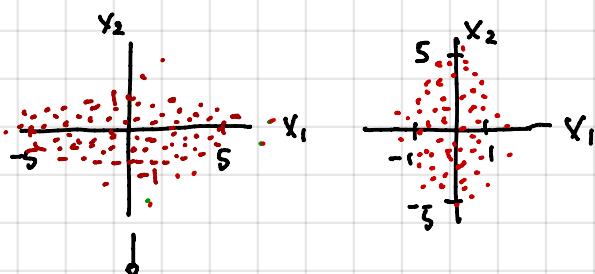
Examples

$$\Sigma_A = \begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma_B = \begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix}$$

$$\Sigma_C = \begin{bmatrix} 5 & 4 \\ 4 & 6 \end{bmatrix}$$

$$\Sigma_D = \begin{bmatrix} 5 & -4 \\ -4 & 6 \end{bmatrix}$$



$\Sigma_{11} = 5$  meaning:

$$\begin{aligned} \text{cov}(x_1, x_1) &= E[(x_1 - \mu_1)^2] \\ &= E[x_1^2] = 5 \end{aligned}$$

The above points are supposedly obtained by random sampling from a Gaussian distribution with  $\bar{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  and  $\Sigma$ .



The definition of eigenvectors:

$$\sum \bar{v} = \lambda \bar{v}$$

$\bar{v}$  = an eigenvector  
 $\lambda$  = an eigen value

→ How can we draw the above figures?  
How can we justify those point distributions?

Q: How to compute  $\lambda$  and  $\bar{v}$  from a given  $\Sigma$ ?

$$A\bar{x} = 0$$

$$\begin{cases} 5x_1 + 6x_2 = 0 \\ 2x_1 + 3x_2 = 0 \end{cases}$$

$$\begin{bmatrix} 5 & 6 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$A\bar{x} = 0 \rightarrow \text{For this:}$$

$\bar{x}$  must be zero!

→ Why? Try to

solve the equations:

$$\begin{cases} 5x_1 + 6x_2 = 0 \\ 2x_1 + 3x_2 = 0 \end{cases}$$

↓

$$\det(\Sigma - \lambda \mathbb{I}) = 0$$

$$\det \begin{pmatrix} \Sigma_{11} - \lambda & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} - \lambda \end{pmatrix} = 0$$

$$\begin{cases} 4x_1 + 6x_2 = 0 \\ 2x_1 + 3x_2 = 0 \end{cases}$$

$$\begin{bmatrix} 4 & 6 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0 \rightarrow \text{even } \bar{x} \neq 0$$

the equation holds

Hence:  $\det(A) = 0$

the ONLY solution  
is  $x_1 = x_2 = 0$ !

$$\begin{aligned} x_1 &= -6/5 x_2 \\ -12x_2 + 3x_2 &= 0 \rightarrow x_2 = 0 \rightarrow x_1 = 0 \end{aligned}$$

Examples:

$$\textcircled{1} \quad \Sigma = \begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow Q: \text{How do we know the distribution directions (or principal axes)?}$$

Answer:

$$\det \begin{pmatrix} 5-\lambda & 0 \\ 0 & 1-\lambda \end{pmatrix} = 0 \rightarrow (5-\lambda)(1-\lambda) = 0$$

$$\lambda_1 = 5$$

$$\lambda_2 = 1$$

$$\text{For } \bar{v}_1 : \sum \bar{v}_i = \lambda_1 \bar{v}_1$$

$$\begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix} \bar{v}_1 = 5 \bar{v}_1 \rightarrow \begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} v_1^a \\ v_1^b \end{bmatrix} = \begin{bmatrix} 5v_1^a \\ 5v_1^b \end{bmatrix}$$

$$\text{Thus: } \bar{v}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \leftarrow$$

$$\begin{bmatrix} 5v_1^a \\ v_1^b \end{bmatrix} = \begin{bmatrix} 5v_1^a \\ 5v_1^b \end{bmatrix} \rightarrow$$

$$v_1^a = 1 \\ v_1^b = 0$$

$$\text{For } \bar{v}_2 : \sum \bar{v}_2 = \lambda_2 \bar{v}_2$$

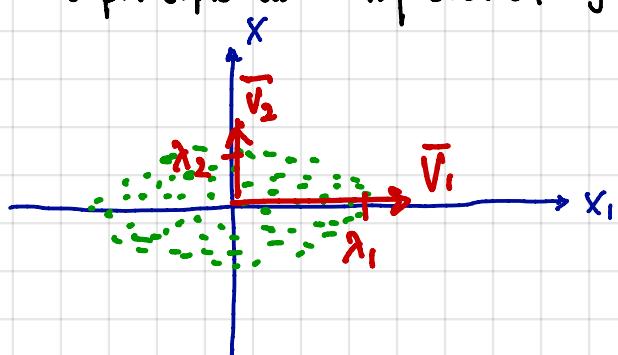
$$\begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} v_2^a \\ v_2^b \end{bmatrix} = \begin{bmatrix} v_2^a \\ v_2^b \end{bmatrix} \rightarrow \begin{bmatrix} 5v_2^a \\ v_2^b \end{bmatrix} = \begin{bmatrix} v_2^a \\ v_2^b \end{bmatrix} \rightarrow$$

$$v_2^a = 0 \\ v_2^b = 1$$

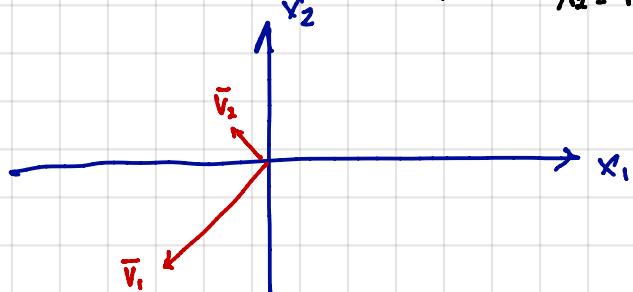
$$\text{Thus: } \bar{v}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



Hence now we know the principle axes represented by  $\bar{v}_1$  &  $\bar{v}_2$ :



$$\textcircled{2} \quad \Sigma = \begin{bmatrix} 5 & 4 \\ 4 & 6 \end{bmatrix} \rightarrow \text{Using the same technique: } \lambda_1 = 9.53 \quad \lambda_2 = 1.47 \rightarrow \bar{v}_1 = (-0.66, -0.75)^T \\ \bar{v}_2 = (-0.75, 0.66)^T$$



# CONDITIONAL & MARGINAL GAUSSIAN DIST.'S

Conditional probability distribution:  $p(x|d)$

$$\text{where: } p(x|d) = \frac{p(x,d)}{p(d)} = \frac{p(d|x)p(x)}{p(d)}$$

Marginal distribution:  $p(d)$

$$\text{where: } p(d) = \int p(x,d) dx = \sum_x p(x,d)$$

Motivation: we know that solving  $p(d)$  and thus  $p(x|d)$  is analytically infeasible & computationally prohibitively expensive.

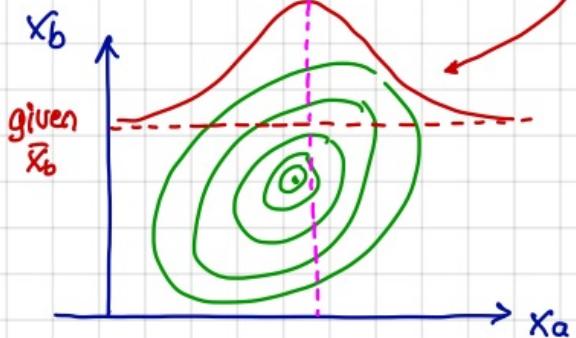


Solution: to use the Gaussian distribution function

## Two Cases

case 1:

We want to know  
the mean & the covariance matrix



We want to know:

1.  $p(x_1 | x_2)$
2.  $p(x_2)$

case 2:

$$\bar{y} = A\bar{x} + \bar{b}$$

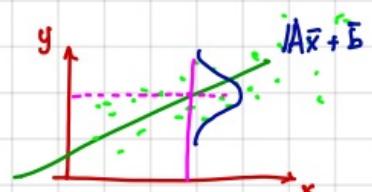
$$p(\bar{x}) = G(\bar{x}; \bar{\mu}, \Lambda^{-1})$$

$$p(\bar{y}|\bar{x}) = G(\bar{y}; A\bar{x} + \bar{b}, L')$$



We want to know:

1.  $p(\bar{x}|\bar{y})$
2.  $p(\bar{y})$



CASE 1 :  $\bar{x} = \begin{pmatrix} \bar{x}_a \\ \bar{x}_b \end{pmatrix}$  and  $p(\bar{x}) = G(\bar{x}; \mu, \Sigma)$   
 $= G\left(\begin{bmatrix} \bar{x}_a \\ \bar{x}_b \end{bmatrix}; \begin{bmatrix} \bar{\mu}_a \\ \bar{\mu}_b \end{bmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}\right)$

Marginal Gaussian distribution:

$$p(\bar{x}_a) = G(\bar{x}_a; \mu_a, \Sigma_{aa})$$

$$p(\bar{x}_b) = G(\bar{x}_b; \mu_b, \Sigma_{bb})$$

Conditional Gaussian distribution:

$$p(\bar{x}_a | \bar{x}_b) = G(\bar{x}_a; \mu_{a|b}, \Sigma_{a|b})$$

$$\text{where: } \mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\bar{x}_b - \mu_b)$$

$$\Sigma_{a|b} = (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1}$$

$$p(\bar{x}_b | \bar{x}_a) = G(\bar{x}_b; \mu_{b|a}, \Sigma_{b|a})$$

$$\text{where: } \mu_{b|a} = \mu_b + \Sigma_{ba} \Sigma_{aa}^{-1} (\bar{x}_a - \mu_a)$$

$$\Sigma_{b|a} = (\Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab})^{-1}$$

CASE 2 :  $\bar{y} = A\bar{x} + \bar{b}$

where:

$$p(\bar{x}) = G(\bar{x}; \bar{\mu}, \Lambda^{-1})$$

$$p(\bar{y} | \bar{x}) = G(\bar{y}; A\bar{x} + \bar{b}, L^{-1})$$

→ Prior

→ Likelihood

Marginal Gaussian distribution:

$$p(\bar{y}) = G(\bar{y}; A\bar{\mu} + \bar{b}, L^{-1} + A\Lambda^{-1}A^T)$$

Conditional Gaussian distribution:

$$p(\bar{x} | \bar{y}) = G(\bar{x}; \bar{\mu}_{x|y}, \Sigma_{x|y})$$

where:

$$\bar{\mu}_{x|y} = (\Lambda + A^T L A)^{-1} (A^T L (\bar{y} - \bar{b}) + \Lambda \bar{\mu})$$

$$\Sigma_{x|y} = (\Lambda + A^T L A)^{-1}$$

# FULL BAYESIAN INFERENCE FOR REGRESSION

In the training stage, we need to estimate:  $\bar{w}$

According to Bayes' rule:

$$p(\bar{w}, \bar{x} | \bar{E}) = \frac{p(\bar{E} | \bar{w}, \bar{x}) p(\bar{w} | \bar{x}) p(\bar{x})}{p(\bar{E})} ; \text{ We don't have any prior information on } \bar{x}$$

$$= \frac{\text{Likelihood}}{\text{Prior}} = \frac{p(\bar{E} | \bar{w})}{p(\bar{w})} p(\bar{w})$$

We know:

- Likelihood:  $p(\bar{E} | \bar{w}) = G(\bar{E}; \phi \bar{w}, \beta^{-1} I)$
- Prior:  $p(\bar{w}) = G(\bar{w}; \alpha, \kappa^{-1} I)$

} we want to know  $p(\bar{w}, \bar{x} | \bar{E})$  or  $p(\bar{w} | \bar{E})$

This problem fits well to CASE 2:

CASE 2:  $\bar{y} = A\bar{x} + \bar{b}$

where:

$$p(\bar{x}) = G(\bar{x}; \bar{\mu}, \Lambda^{-1}) \rightarrow \text{Prior}$$

$$p(\bar{y} | \bar{x}) = G(\bar{y}; A\bar{x} + \bar{b}, L^{-1}) \rightarrow \text{Likelihood}$$

Marginal Gaussian distribution:

$$p(\bar{y}) = G(\bar{y}; A\bar{\mu} + \bar{b}, L^{-1} + A\Lambda^{-1}A^T)$$

Conditional Gaussian distribution:

$$p(\bar{x} | \bar{y}) = G(\bar{x}; \bar{\mu}_{x|y}, \Sigma_{x|y})$$

where:

$$\bar{\mu}_{x|y} = (\Lambda + A^T L A)^{-1} (A^T L (\bar{y} - \bar{b}) + \Lambda \bar{\mu})$$

$$\Sigma_{x|y} = (\Lambda + A^T L A)^{-1}$$

Where:

$$\mu \rightarrow 0$$

$$\Lambda^{-1} \rightarrow \alpha^{-1} I$$

$$A \rightarrow \phi$$

$$\bar{x} \rightarrow \bar{w}$$

$$\bar{y} \rightarrow \bar{E}$$

$$\bar{b} \rightarrow 0$$

$$L \rightarrow \beta I$$

Hence:  $p(\bar{w} | \bar{E}) = G(\bar{w}; \bar{m}_N, \mathbb{S}_N)$

where:  $\mathbb{S}_N = (\alpha I + \phi^T \beta I \phi)^{-1} = (\alpha I + \beta \phi^T \phi)^{-1}$

$$\bar{m}_N = \mathbb{S}_N (\phi^T \beta I \phi + \alpha I \phi) = \mathbb{S}_N \beta \phi^T \bar{E}$$

Hence:

$$\bar{m}_N = \beta S_N \Phi^T \bar{e}$$

Mx1 MxM MXN Nx1

$$S_N = (\alpha I + \beta \Phi^T \Phi)^{-1}$$

MxM MxM MxM

In other words:

$$\bar{m}_N = \beta (\alpha I + \beta \Phi^T \Phi)^{-1} \Phi^T \bar{e}$$

MLE:

$$\bar{w}_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T \bar{e}$$

$$\bar{w}_{MAP} = \beta (\alpha I + \beta \Phi^T \Phi)^{-1} \Phi^T \bar{e}$$

they are the same!

Note: Unlike MAP, full Bayesian provides  $S_N$  which in turn gives us:  $p(\bar{w} | \bar{e})$ . This means we can obtain uncertainty.

# ● SEQUENTIAL LEARNING

## [Tr] Sequential Bayesian Learning

Concept:

simplification {

$$\begin{aligned} & p(\bar{w} | t_1, t_2, t_3, x_1, x_2, x_3) \rightarrow \text{we want to estimate } \bar{w}. \text{ the curve parameters.} \\ & = p(w | t_1, t_2, t_3) \\ & = \frac{p(t_3 | t_2, t, \bar{w}) p(t_2, t_1, \bar{w})}{p(t_3, t_2, t_1)} \\ & = \frac{p(t_3 | \bar{w})}{p(t_3)} \frac{p(t_2 | t_1, \bar{w})}{p(t_2)} p(t_1, \bar{w}) \\ & = \frac{p(t_3 | \bar{w})}{p(t_3)} \frac{p(t_2 | \bar{w})}{p(t_2)} \frac{p(t_1 | \bar{w})}{p(t_1)} p(\bar{w}) \end{aligned}$$

For illustration, see the textbook page #155.

Q: What is the implication of the previous formulation?

A:  $p(\bar{w} | t_1, t_2, t_3) = \frac{p(t_3 | \bar{w})}{p(t_3)} \frac{p(t_2 | \bar{w})}{p(t_2)} \frac{p(t_1 | \bar{w})}{p(t_1)} p(\bar{w})$

↑ prior

↑ likelihood 1  
evidence 1

↑ likelihood 2  
evidence 2

↑ likelihood 3  
evidence 3

posterior 1:  
 $p(\bar{w} | t_1) = \frac{p(t_1 | \bar{w}) p(\bar{w})}{p(t_1)}$

posterior 2:  
 $p(\bar{w} | t_2, t_1) = \frac{p(t_2 | \bar{w}) p(\bar{w} | t_1)}{p(t_2)}$

posterior 3:  
 $p(\bar{w} | t_3, t_2, t_1) = \frac{p(t_3 | \bar{w}) p(\bar{w} | t_2, t_1)}{p(t_3)}$

All these imply sequential learning or online learning, since we can process once we have a new data, and update the old stored value of  $\bar{m}_N$  and  $S_N$ .

Q : Practically how can we implement it?

A : Like the discussion above, we need to define :

① Gaussian prior :

$$p(\bar{w}) = G(\bar{w}; \bar{m}_0, \alpha^{-1} \mathbb{I});$$

$\downarrow \quad \downarrow$   
 $\mu = \bar{m}_0 \quad \Lambda = \alpha \mathbb{I}$

② Gaussian likelihood :

$$p(t_n | \bar{w}) = G(t_n; \bar{w}^T \bar{\phi}(x_n), \beta^{-1})$$

Once  $t_1$  arrives, we compute :

$$p(\bar{w} | t_1) = G(\bar{w}; \bar{m}_1, \mathbb{S}_1)$$

See the definitions  
of  $\bar{m}_1$  &  $\mathbb{S}_1$   
in the previous note.

Once  $t_2$  arrives, we compute :

$$p(\bar{w} | t_2) = G(\bar{w}; \bar{m}_2, \mathbb{S}_2)$$

But, what are  $\bar{m}_2, \mathbb{S}_2$  here?

Q:  $p(\bar{w} | t_2) = G(\bar{w}; \bar{m}_2, \mathbb{S}_2) \rightarrow$  how to compute this?

A : We know:  $p(w | t_2) = \frac{p(t_2 | w) p(\bar{w} | t_1)}{p(t_2)}$

$$\text{where: } p(\bar{w} | t_1) = G(\bar{w}; \bar{m}_1, \mathbb{S}_1); \quad p(t_2 | w) = G(t_2; \bar{w}^T \bar{\phi}(x_2), \beta^{-1})$$

case 2

Hence:



$$\Phi_2 = \begin{bmatrix} \bar{\phi}^T(x_2) \\ \vdots \\ \bar{\phi}^T(x_2) \end{bmatrix}_{M \times 1}$$

$$p(\bar{w} | t_2) = G(\bar{w}; \bar{m}_2, \mathbb{S}_2)$$

with:

$$\left\{ \begin{array}{l} \mathbb{S}_2^{-1} = \mathbb{S}_1^{-1} + \beta \Phi_2^T \Phi_2 \\ M \times M \qquad M \times M \end{array} \right. \quad \left. \begin{array}{l} \bar{m}_2 = \mathbb{S}_2 (\mathbb{S}_1^{-1} \bar{m}_1 + \beta \Phi_2^T t_2) \\ M \times M \qquad M \times 1 \end{array} \right.$$

Q: How can we transform the above formulation to an algorithm?

A : We simply use the definitions of  $\mathbb{S}_2^{-1}$  and  $\bar{m}_2$  above sequentially:

$$\mathbb{S}_n^{-1} = \mathbb{S}_{n-1}^{-1} + \beta \Phi_n^T \Phi_n$$

$$\bar{m}_n = \mathbb{S}_n (\mathbb{S}_n^{-1} \bar{m}_{n-1} + \beta \Phi_n^T t_n)$$

Note :  $p(\bar{w}|t_1) = \underbrace{G(\bar{w}; \bar{M}_1, S_1)}_{\text{Gaussian prior}} ; p(t_2|w) = \underbrace{G(t_2; \bar{w}^T \phi(x_2), \beta^{-1})}_{\text{Gaussian likelihood}}$

CASE 2 :  $\bar{y} = A\bar{x} + \bar{b}$

where :  $p(\bar{x}) = G(\bar{x}; \bar{\mu}, \Lambda^{-1}) \rightarrow \text{Prior}$   
 $p(\bar{y}|\bar{x}) = G(\bar{y}; A\bar{x} + \bar{b}, L^{-1}) \rightarrow \text{Likelihood}$

Marginal Gaussian distribution :

$$p(\bar{y}) = G(\bar{y}; A\bar{\mu} + \bar{b}, L^{-1} + A\Lambda^{-1}A^T)$$

Conditional Gaussian distribution :

$$p(\bar{x}|\bar{y}) = G(\bar{x}; \bar{\mu}_{x|\bar{y}}, \Sigma_{x|\bar{y}})$$

where :  $\bar{\mu}_{x|\bar{y}} = (\Lambda + A^T L A)^{-1} (A^T L (\bar{y} - \bar{b}) + \Lambda \bar{\mu})$   
 $\Sigma_{x|\bar{y}} = (\Lambda + A^T L A)^{-1}$

$$\begin{aligned} \bar{x} &\rightarrow \bar{w} \\ \bar{y} &\rightarrow t_2 \\ \Lambda &\rightarrow S_1^{-1} \\ A &\rightarrow \Phi(\bar{x}_2) \rightarrow \Phi_2 \\ &\quad M \times 1 \qquad N \times M ; N=1 \\ L &\rightarrow \beta \\ b &\rightarrow 0 \\ M &\rightarrow M_1 \end{aligned}$$

Hence : Posterior :

$$p(\bar{w}|t_2) = G(\bar{w}; \bar{M}_2, S_2)$$

where :

$$S_2^{-1} = \left( S_1^{-1} + \Phi_2^T \beta \Phi_2 \right)^{-1}$$

$$\bar{M}_2 = S_2^{-1} (\Phi_2^T \beta t_2)$$

# PREDICTIVE DISTRIBUTION

Motivation: Instead of  $\bar{w}$  and its probability,  $p(\bar{w} | \bar{E})$ , we are more interested in  $t_{\text{new}}$  and its probability,  $p(t_{\text{new}} | \bar{E})$ . This  $p(t_{\text{new}} | \bar{E})$  is not obtained in the previous discussion.

$$t_* = t_{\text{new}} = Y_* + \varepsilon = \bar{w}^T \bar{X}_* + \varepsilon :$$

$$p(t_* | X_*, \bar{E}, \bar{x}, \alpha, \beta) = \int p(t_* | \bar{w}) = \int p(t_*, \bar{w} | \bar{E}) d\bar{w}$$

Marginalization over all possible values of  $\bar{w}$ .

$$= \int \underbrace{p(t_* | \bar{w})}_{\textcircled{1}} \underbrace{p(\bar{w} | \bar{E})}_{\textcircled{2}} d\bar{w}$$

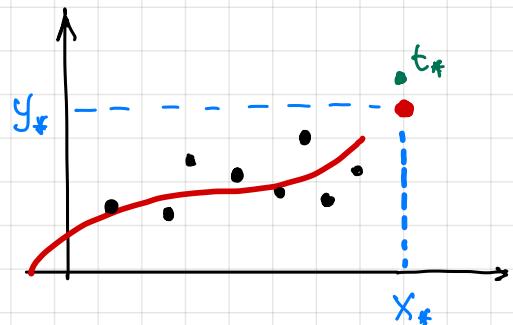
$$\textcircled{1} \quad p(t_* | \bar{w}) = G(t_*; \bar{w}^T \Phi(x_*), \beta^{-1})$$

$$\textcircled{2} \quad p(\bar{w} | \bar{E}) = \frac{p(\bar{E} | \bar{w}) p(\bar{w})}{p(\bar{E})}$$

where:

$$p(\bar{w}) = G(\bar{w}; \bar{m}_0, \bar{S}_0^{-1})$$

$$p(\bar{E} | \bar{w}) = G(\bar{E}; \bar{\Phi} \bar{w}, \beta^{-1} \bar{I})$$



According to Bayesian sequential learning:

$$p(\bar{w} | \bar{E}, \alpha, \beta) = G(\bar{w}; \bar{m}_n, \bar{S}_n)$$

where:

$$\bar{S}_n^{-1} = \bar{S}_{n-1}^{-1} + \beta \bar{\Phi}^T \bar{\Phi}$$

$$\bar{m}_n = \bar{S}_n (\bar{S}_{n-1}^{-1} \bar{m}_{n-1} + \beta \bar{\Phi}^T \bar{t})$$

$$\left. \begin{aligned} \bar{m}_0 &= 0 \\ \bar{S}_0^{-1} &= \alpha \end{aligned} \right\}$$

If the prior and likelihood are Gaussians:

$$p(t_* | \bar{t}) = \int p(t_* | \bar{w}) p(\bar{w} | \bar{t}) d\bar{w}$$

① Gaussian Likelihood (w.r.t.  $\bar{w}$ )      ② Gaussian Prior (w.r.t.  $\bar{w}$ )

then, according to case 2, the marginal distribution is:

case 2 :

Recall the definition of the Marginal Gaussian distribution (lecture note 7):

Given:  $p(x) = G(x; \mu, \Lambda^{-1})$

$p(y|x) = G(y; Ax + b, L^{-1})$

Then:  $p(y) = G(y; A\mu + b, \bar{\Lambda}^{-1} + A\Lambda^{-1}A^T)$

$p(t_* | x_*, \bar{t}, \alpha, \beta) = G(t_*; \bar{m}_*, \bar{s}_*)$

$\bar{m}_* = \bar{m}_N^T \bar{\phi}(x_*)$

$\bar{s}_* = \beta^{-1} + \bar{\phi}^T(x_*) \mathbb{S}_N \bar{\phi}(x_*)$

There are 2 steps: ① Prediction  
② Update

See textbook page #156.



This implies that given the training data  $\{(x_i, t_i) \dots (x_n, t_n)\}$   
we can interpolate or extrapolate the prediction of  $t$  and  
know the uncertainty of the prediction!

see textbook figure 3.8 (page #157)

If there are two data  $[t_{*1}, t_{*2}]$  we want to process at once:

$$\bar{m}_* = \Phi_{1,2} \bar{m}_N \quad ; \quad \bar{s}_* = \beta^{-1} \mathbb{I} + \Phi_{1,2}^T \mathbb{S}_N \Phi_{1,2}$$

$2 \times 1 \quad 2 \times M \quad M \times 1 \quad 2 \times 2 \quad 2 \times M \quad M \times M \quad M \times 2$

Note:

Q: Why is  $p(t_* | \bar{E})$  beneficial?

A: If we know  $p(t_* | \bar{E}) = G(t_*; \bar{m}_*, s_*)$ , it means we obtain:



This means: we know the probability of all possible prices, and we know the most probable price.

→ (Prediction & its uncertainty)

---

Q: fhw can we know that:  $p(t_* | \bar{E}) = \int p(t_*, \bar{w} | \bar{E}) d\bar{w}$ ?

A: Basic idea of  $p(t_* | \bar{E})$  is that we need to find the correlation between  $t_*$  and  $\bar{E}$ .

We know that in the testing, to estimate  $t_*$  we need  $x_*$  and  $\bar{w}$  to obtain  $t_*$ ; Also, in the training, we can get  $\bar{w}$  from  $\bar{E}$ .

Hence, the correlation of  $t_*$  &  $\bar{E}$  can be obtained through  $\bar{w}$ .

Yet, at the same time we don't care about  $\bar{w}$  in the end.

Therefore, we marginalize over  $\bar{w}$  to obtain  $p(t_* | \bar{E})$ .